



9-27-00

#

REFERENCE & ASSOCIATES
PATENT FILING TRANSMITTAL

DOCKET NO. YOR920000390
(590.023)

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

Box Patent Application
Commissioner of Patents and Trademarks
Washington, D.C. 20231



PATENT FILING TRANSMITTAL

Transmitted herewith for filing is the Patent Application of: Mukund Padmanabham, George A. Saon,
and Geoffrey G. Zweig

For: LATTICE-BASED UNSUPERVISED MAXIMUM LIKELIHOOD LINEAR REGRESSION
FOR SPEAKER ADAPTATION

TYPE OF FILING

This new patent application is for a(n):

- ☒ Utility
- ☐ Design
- ☐ Plant
- ☐ Divisional
- ☐ Continuation
- ☐ Continuation-in-part

Benefit of a prior filed application

- ☐ This application claims the benefit of an earlier filed U.S. Patent Application under 35 USC 120.
- ☐ Please accord Applicant the benefit of the priority date of _____ to this case pursuant to 35 USC 119. Applicant's claim for priority is based on application _____ filed in _____ on that date.

Filing under 37 CFR 1.53 (Utility) or 37 CFR 1.153 (Design)

- ☒ This is an application filed pursuant to 37 CFR 1.53 or 37 CFR 1.153, permitting receipt of a filing date upon filing of a specification, at least one claim and necessary drawings.
- ☒ In the event any parts of this application are incomplete, please treat this as a filing under 37 CFR 1.53 or 37 CFR 1.153.

ENCLOSURES

- ☒ 20 - pages of written description;
- ☒ 5 - pages of claims;
- ☒ 1 - pages of abstract;
- ☐ _____ - sheets of formal drawings;
- ☒ 2 - sheets of informal drawings;
- ☒ Declaration and Power of Attorney or listing of inventors;
- and**
- ☒ Two postcards for return to us as proof of receipt of the above documents.

plus

- ☒ An Assignment of the invention to IBM Corporation and an Assignment cover sheet;

- ☐ Verified Statement Claiming Small Entity Status (37 CFR 1.9(f) and 1.27(b))
☐ Form PTO-1449 (IDS) and two copies of the references listed thereon;
☐ A certified copy of _____ (country) patent application number _____ (priority document).
☐ A preliminary amendment;
☐ Declaration of Biological Deposit;
☐ Submission of sequence listing, computer readable copy and/or amendment relating thereto for biotechnology invention containing nucleotide and/or amino acid sequence;
☐ An associate power of attorney;
☐ Other.

DECLARATION OR OATH

The enclosed Declaration or Oath has been executed by:

- ☒ Inventor(s);
☐ Legal representative of the inventors (37 CFR 1.42 or 1.43);
☐ Joint inventor or person showing proprietary interest on behalf of an inventor who refused to sign or who cannot be reached and this is a petition required by 37 CFR 1.47 and the statement required by 37 CFR 1.47 is attached;
☐ Has not been executed and is enclosed for the purposes of identifying the inventors.

INVENTORSHIP STATEMENT

The inventorship for all the claims in this application is:

- ☒ the same;
☐ not the same and, as an explanation, a statement is/ will be submitted.

LANGUAGE

The application submitted herewith is:

- ☒ in English;
☐ in not in English and in terms of 37 CFR 1.52(d) a verified translation is
☐ attached
☐ not attached.

FEE CALCULATION

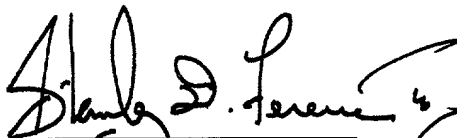
The filing fee has been calculated as shown below:

				SMALL ENTITY OR		OTHER THAN A SMALL ENTITY	
				RATE	FEE	RATE	FEE
BASIC FEE Design Patent				\$155	\$	\$310	\$
BASIC FEE Utility Patent				\$345	\$	\$690	\$690
EXTRA FEES				RATE	FEE	RATE	FEE
TOTAL CLAIMS	15	MINUS 20=	0	x 9=	\$0	x18=	\$
INDEP.CLAIMS	3	MINUS 3 =	0	x 39=	\$0	x78=	\$
<input type="checkbox"/> MULTIPLE DEP.CLAIM				+135=	\$	+270=	\$
<input checked="" type="checkbox"/> ASSIGNMENT				+ 40=	\$	+40=	\$
<input type="checkbox"/> RULE 53 SURCHARGE				+ 65=	\$	+130=	\$
TOTAL					\$		\$730

FEE PAYMENT

- ☐ Attached is Check No. _____ in the sum of \$ _____ to cover the filing fee and, if applicable, the assignment fee.
- ☒ Please charge **IBM Corporation Deposit Account No. 50-0510** the amount of \$730.00, to cover the filing fee. Duplicate copies of this letter are enclosed. In the event of non-payment or improper payment of a required fee, the Commissioner is authorized to charge or credit **Deposit Account No. 50-0510** as required to correct the error.

Respectfully submitted,



Stanley D. Ference III
Reg. No. 33,879

Dated: September 26, 2000

FERENCE & ASSOCIATES
129 Oakhurst Road
Pittsburgh, Pennsylvania 15215
(412) 781-7386
(412) 781-8390-Facsimile

PATENT

Docket No. YOR920000390US1
(590.023)

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

Applicant(s) : Mukund Padmanabhan et al Group Art: not yet assigned
Serial No. : not yet assigned Examiner: not yet assigned
Filed : herewith
For : LATTICE-BASED UNSUPERVISED MAXIMUM LIKELIHOOD
LINEAR REGRESSION FOR SPEAKER ADAPTATION

EXPRESS MAIL CERTIFICATE

Express Mail Label No. EL503717394US

Date of Deposit 26 September 2000

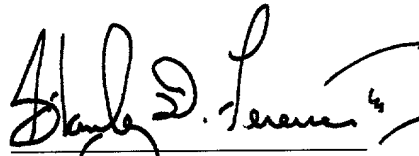
I hereby certify that the following attached paper(s) or fee:

Patent Application
Written Description
Claims 1-5
Abstract
Drawings (Figs. 1-4)
Declaration and Power of Attorney
Assignment
Patent Filing Transmittal
Certificate of Express Mail
Two Return Postcards

are being deposited with the United States Postal Service "Express Mail Post Office to Addressee" service under 37 C.F.R. 1.10 on the date indicated above and is addressed to the Assistant Commissioner for Patents, Washington, D.C. 20231.

Stanley D. Ference III

(Typed or printed name of
person mailing paper)



(Signature of person mailing
paper(s) or fee)

Mailing Address:

FERENCE & ASSOCIATES
129 Oakhurst Road
Pittsburgh, Pennsylvania 15215
(412) 781-7386
(412) 781-8390-Facsimile

LATTICE-BASED UNSUPERVISED MAXIMUM LIKELIHOOD LINEAR

REGRESSION FOR SPEAKER ADAPTATION

Field of the Invention

The present invention generally relates to methods and arrangements for providing
5 speaker adaptation in connection with speech recognition.

Background of the Invention

Acoustic adaptation is playing an increasingly important role in speech recognition
systems, to compensate for the acoustic mismatch between training and test data, and also
to adapt speaker-independent systems to individual speakers. Most speech recognition
10 systems use acoustic models that include multi-dimensional gaussians that model the
probability density function (pdf) of the feature vectors for different classes. (For general
background on speech recognition, including gaussian mixture pdf's, see, e.g.
Fundamentals of Speech Recognition, Lawrence Rabiner and Biing-Hwang Juang,
Prentice Hall, 1993; and *Statistical Methods for Speech Recognition*, Frederick Jelinek,
15 The MIT Press, 1997.) A commonly used adaptation technique in this connection is
maximum likelihood linear regression (MLLR), which assumes that the parameters of the
gaussians are transformed by an affine transform into parameters that better match the test

or adaptation data. In a simple implementation, the mean u_i of each gaussian g_i is transformed according to $u_i' = Au_i$ where A is the transform matrix, and u_i is optionally padded with ones to represent an offset. The transform is chosen so as to maximize the probability of a collection of adaptation data with associated transcriptions. In more

5 sophisticated implementations, the gaussian variances may also be adjusted. MLLR is further discussed, for instance, in Leggetter et al., "Speaker Adaptation of Continuous Density HMM's Using Multivariate Linear Regression", Proceedings of ICSLP '94, Yokohama, Japan, 1994. This technique is also often used in "unsupervised" mode, where the correct transcription of the adaptation data is not known, and a first pass decoding

10 using a speaker independent system is used to produce an initial transcription.

Although MLLR appears to work fairly well even when the unsupervised transcription is mildly erroneous, it is recognized herein that further improvements are possible.

Accordingly, a need has been recognized, *inter alia*, in connection with improving

15 upon the shortcomings and disadvantages associated with conventional arrangements such as those discussed above.

Summary of the Invention

In accordance with at least one presently preferred embodiment of the present invention, it is presently recognized that it is possible to improve upon the performance of MLLR, even when the unsupervised transcription is mildly erroneous, by taking into
5 account the fact that the initial transcription contains errors. This may be accomplished, for example, by considering not just the “1-best” (*i.e.*, single best) transcription produced during the first pass decoding, but the top N candidates. (See, for example, Jelinek [1997], *supra*, for a description of “N-best” decoding and definitions associated therewith.) Alternatively, if the first pass decoding produces a word graph, this can be
10 used as the reference word graph, instead of the 1-best or N-best reference transcriptions. In contrast to an N-best list, which simply enumerates a relatively small number (e.g. 100 or 1000) of likely word sequences, a word graph is a compact representation of all the word sequences that have any appreciable probability. An example of a word graph is illustrated in Figure 2.

15 Broadly contemplated herein, in accordance with at least one presently preferred embodiment of the present invention, is a formulation that affinely transforms the means of the gaussians to maximize the log likelihood of the adaptation data under the assumption that a word graph is available that represents all possible word sequences that correspond

to the adaptation data. The word graph is produced during a first pass decoding with speaker independent models. It is also possible to consider only those regions of the word graph that represent a high confidence of being correct to further improve the performance.

5 In one aspect, the present invention provides a method of providing speaker adaptation in speech recognition, the method comprising the steps of: providing at least one speech recognition model; accepting speaker data; generating a word lattice based on the speaker data; and adapting at least one of the speaker data and the at least one speech recognition model in a manner to maximize the likelihood of the speaker data with respect
10 to the generated word lattice.

 In another aspect, the present invention provides an apparatus for providing speaker adaptation in speech recognition, the apparatus comprising: at least one speech recognition model; an accepting arrangement which accepts speaker data; a lattice generator which generates a word lattice based on the speaker data; and a processing
15 arrangement which adapts at least one of the speaker data and the at least one speech recognition model in a manner to maximize the likelihood of the speaker data with respect to the generated word lattice.

Furthermore, in another aspect, the present invention provides a program storage device readable by machine, tangibly embodying a program of instructions executable by the machine to perform method steps for providing speaker adaptation in speech recognition, the method comprising the steps of: providing at least one speech recognition
5 model; accepting speaker data; generating a word lattice based on the speaker data; and adapting at least one of the speaker data and the at least one speech recognition model in a manner to maximize the likelihood of the speaker data with respect to the generated word lattice.

For a better understanding of the present invention, together with other and further
10 features and advantages thereof, reference is made to the following description, taken in conjunction with the accompanying drawings, and the scope of the invention will be pointed out in the appended claims.

Brief Description of the Drawings

Fig. 1 illustrates a Hidden Markov Model (HMM) structure used to generate
15 Maximum A-Posteriori Probability (MAP) lattices;

Fig. 2 illustrates word traces produced by the MAP lattice HMM, and their connection into a word lattice;

Fig. 3 illustrates a histogram of state posterior probabilities; and

Fig. 4 illustrates a graph of word error rate versus confidence threshold.

Description of the Preferred Embodiments

Throughout the present disclosure, various terms are utilized that are generally well-known to those of ordinary skill in the art. For a more in-depth definition of such terms, any of several sources may be relied upon, including Rabiner and Juang (1993), *supra*, and Jelinek (1997), *supra*. First, a theoretical framework is discussed in connection with at least one presently preferred embodiment of the present invention.

Also, the article "Lattice Based Unsupervised MLLR for Speaker Adaptation in Speech Recognition Systems" (Mukund Padmanabhan, George Saon, Geoffrey Zweig, ISCA ITRW ASR2000 Paris, France 2000 [<http://www-tlp.limsi.fr/asr2000>]) is hereby fully incorporated by reference as if set forth in its entirety herein.

In a typical speech recognition system, the speech signal is represented as a sequence of observation vectors, and throughout the present discussion, y_t denotes the multi-dimensional observation at time t , and y_1^T denotes the T observations corresponding to the adaptation data. The pdf's of each context dependent phonetic state s is modeled by a single gaussian (this can be easily generalized to mixtures of gaussians) with mean

and diagonal covariance μ_s , Λ_s . θ is used to indicated the current values of the gaussian parameters, and $\hat{\theta}$ is used to denote the future (adapted) values to be estimated. The probability density of the observation y_t given the pdf of state s is denoted $p_\theta(y_t / s)$. It will presently be assumed that θ and $\hat{\theta}$ are related in the following way: $\hat{\mu}_s = A\mu_s$, $\hat{\Lambda}_s =$
5 Λ_s , i.e., only the current means of the gaussians are linearly transformed, and all means are transformed by the same matrix A.

Typically, in an MLLR framework, the general objective is defined as follows:
given a transcription w of the adaptation data, find $\hat{\theta}$ (or equivalently A) so that the log likelihood of the adaptation data, y_I^T is maximized. The transcription w can be represented
10 as a sequence of K states $s_1 \dots s_K$, and the T observation frames can be aligned with this sequence of states. However, the alignment of the T frames with the sequence of states is not known. Let s_t denote the state at time t . The objective is to find the maximum likelihood transform theta, and can now be written as:

$$\begin{aligned}
\hat{\theta}^* &= \arg \max_{\hat{\theta}} \log \left[p_{\hat{\theta}}(y_1^T) \right] \\
&= \arg \max_{\hat{\theta}} E_{s_1^T / y_1^T, \theta} \log \left[p_{\hat{\theta}}(y_1^T) \right] \\
&= \arg \max_{\hat{\theta}} \sum_{s_1^T} p_{\theta}(s_1^T / y_1^T) \log \left[p_{\hat{\theta}}(y_1^T, s_1^T) \right]
\end{aligned} \tag{1}$$

In a lattice-based MLLR currently contemplated in accordance with at least one embodiment of the present invention, it is assumed that the word sequence, and thus the state sequence s_1^K , corresponding to the adaptation data cannot be uniquely identified, and this uncertainty is incorporated in the form of a lattice or word graph. Preferably, the word graph is produced by a first pass decoding with speaker independent models. The formulation of the maximum likelihood problem is essentially identical to equation (1), but with one significant difference. In (1), the states s_t were assumed to belong to the alphabet of K states $s_1 \dots s_K$, with the only allowed transitions being $s_t \rightarrow s_t$ and $s_t \rightarrow s_{t+1}$. In a lattice-based MLLR formulation according to at least one presently preferred embodiment of the present invention, the transition between the states is dictated by the structure of the word graph. Additionally, it is possible to take into account the language model probabilities (which are ignored in the MLLR formulation), by incorporating them into the transition probability corresponding to the transition from the final state of a word in the word graph to the initial state of the next connected word in the word graph.

The disclosure now turns to a decoding strategy for producing word graphs.

In accordance with at least one presently preferred embodiment of the present invention, a Maximum A-Posteriori Probability (MAP) word lattice is preferably generated using word internal acoustic models and a bigram language model. MAP lattices and

bigram language models are discussed generally in several publications, including (Jelinek, 1997). To construct the lattice, it may be assumed that the utterance in question is produced by an HMM with a structure such as that shown in Figure 1. Each pronunciation variant in the vocabulary appears as a linear sequence of phones in the HMM, and the structure of this model permits the use of word-internal context dependent phones. Preferably, a bigram language model is used with modified Kneser-Ney smoothing (*see*, for example, Kneser and Ney, "Improved Backing-off for n-gram Language Modeling", *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing*, 1995) . Here, there is an arc from the end of each word to a null word-boundary state, and this arc has a transition probability equal to the back-off probability for the word. From the word-boundary state, there is an arc to the beginning of each word, labeled with the unigram probability. For word pairs for which there is a direct bigram probability, an arc is preferably introduced from the end of the first word to the beginning of the second, and this arc preferably has a transition probability equal to the discounted bigram probability. Preferably, the dynamic range of the acoustic and language-model probabilities is normalized by using an appropriate language model weight, such as 15.

Preferably, the MAP lattice is constructed by computing the posterior state occupancy probabilities for each state at each time:

$$P(S_t = s | y_1^T) = \frac{\alpha_s^t \beta_s^t}{P(y_1^T)}$$

where $\alpha_s^t = P(y_1^t, St=s)$ and $\beta_s^t = P(y_{t+1}^T / St=s)$, and then computation posterior word occupancy probabilities by summing over all the states interior to each word. That is, if w_t is the set of states in word W_t , then the following is preferably computed at each time

5 frame:

$$\sum_{s \in w_t} \frac{\alpha_s^t \beta_s^t}{P(y_1^T)}$$

Preferably, the N likeliest words are kept track of at each frame, and these are preferably output as a first step in the processing.

It will be noticed that a word will be on the list of “likeliest words” for a period of

10 time, and thereafter will “fall off” that list. Thus, the output of the first step may preferably be a set of word traces, as illustrated in Figure 2. The horizontal axis is time, while the vertical axis ranges over all the pronunciation variants.

Preferably, the next step will be to connect the word traces into a lattice. Many connection schemes are possible, but it has been found that the following strategy is quite

15 effective. It requires that one more quantity be computed as the word traces are

generated: the temporal midpoint of each trace as computed from the first moment of its posterior probability:

$$\frac{\sum_{t=start}^{t=end} tP_t(W)}{\sum_{t=start}^{t=end} t}$$

To construct an actual lattice, a connection is preferably added from the end of one word trace to the beginning of another if the two overlap, and the midpoint of the second is to the right of the midpoint of the first. This is illustrated at the bottom of Figure 2. (It has also been found to be convenient to discard traces that do not persist for a minimum period of time, or which do not reach an absolute threshold in posterior probability.)

To evaluate the lattices, the oracle worderror rate (i.e. the error rate of the single path through the lattice that has the smallest edit distance from the reference script) can be computed. This is the best worderror rate that can be achieved by any subsequent processing to extract a single path from the lattice. For voicemail transcription, the MAP lattices have an oracle word error rate of about 9%, and the ratio of the number of word occurrences in the lattices to the number of words in the reference scripts is about 64.

Due to the rather lax requirements for adding links between words, the average indegree for a word is 74; that is, there are about 74 possible predecessors for each word in the graph. The MAP lattice that is produced in this way is suitable for a bigram language

model: the arcs between wordends can be labeled with bigram transition probabilities, but is too large for a straightforward expansion to trigram context. In order to reduce its size, a second pass is preferably made, where the posterior probability of transitioning along the arcs that connect wordtraces is computed. That is, if s_t is the last state in one word trace
 5 and s_j is the first state in a successor and a_{ij} is the weighted language model transition probability of seeing the two words in succession, one may compute:

$$P(S_t = s_i, S_{t+1} = s_j | y_1^T) = \frac{\alpha_{s_i}^t \beta_{s_j}^{t+1} a_{ij} b_j(y_{t+1})}{P(y_1^T)}$$

The above equation represents the posterior probability of being in state s_i at time t and in state s_j at time $t + 1$, and transitioning between the words at an intermediate time.
 10 For each link between word traces, this quantity is summed over all time to get the total probability that the two words occurred sequentially; the links with the lowest posteriors are then discarded. It should be noted that a separate quantity is preferably computed for every link in the lattice. Thus, even if two links connect traces with the same word labels, the links will in general receive different posterior probabilities because the traces will lie
 15 in different parts of the lattice, and therefore tend to align to different segments of the acoustic data.

As in Mangu et al., "Lattice Compression in the Consensual Post-Processing Framework" (Proceedings of SCI/ISAS, Orlando, Fla., 1999), it has been found that over 95% of the links can be removed without a major loss of accuracy. Here, it was found that pruned lattices had an average indegree a little under 4, and an oracle error rate of
5 about 11%. After pruning, lattices were expanded to a trigram context, and the posterior state occupancy probabilities needed for MLLR were computed with a modified Kneser-Ney trigram language model, along with leftword context dependent acoustic models.

Accordingly, the disclosure now turns to a discussion of a confidence-related pruning method that enables regions of low confidence to be discarded.

10 Word lattices have been used in a variety of confidence estimation schemes (see, for example, Kemp et. al., "Estimating Confidence Using Word Lattices" (Proceedings of ICASSP '97, 1997) and Evermann et al., "Large Vocabulary Decoding and Confidence Estimation Using Word Posterior Probabilities (Proceedings of ICASSP '00, 2000). Here, the simplest possible measure posterior phone probability was explored for discarding
15 interpretations in which there was low confidence. It is to be recalled that, as a first step in MLLR, the posterior gaussian probabilities are computed for all the gaussians in the system. This is computed on a phone-by-phone basis, first computing the posterior phone

probability, and then multiplying by the relative activations for the gaussians associated with the phone. For phone s_i with gaussian mixture G_i , and for a specific time frame y_t ,

$$P(G_i = g_j | y_1^T) = P(S_i = s_i | y_1^T) \frac{g_j(y_t)}{\sum_{g \in G_i} g(y_t)}$$

Since the gaussian posteriors are used to define a set of linear equations that are
5 solved for the MLLR transform, it is reasonable to assume that noisy or uncertain estimates of the posteriors will lead to a poor estimate of the MLLR transform. To examine the truth of this hypothesis, the MLLR transform was estimated from subsets of the data, using only those estimates of $P(S_i = s_i | y_1^T)$ that were above a threshold, typically 0.7 to 0.9.

10 The experiments were performed on a voicemail transcription task. (For a general discussion of voicemail transcription, see Padmanabhan et al., "Recent Improvements in Voicemail Transcription" (Proceedings of EUROSPEECH '99, Budapest, Hungary, 1999). The speaker independent system has 2313 context dependent phones, and 134,000 diagonal gaussian mixture components, and was trained on approximately 70 hours of
15 data. The feature vectors are obtained in the following way: 24 dimensional cepstral vectors are computed every 10ms (with a window size of 25ms). Every 9 consecutive cepstral vectors are spliced together forming a 216 dimensional vector which is then

projected down to 39 dimensions using heteroscedastic discriminant analysis and maximum likelihood linear transforms (see Saon et al., "Maximum Likelihood Discriminant Feature Spaces", to appear in Proceedings of ICASSP '2000, Istanbul, 2000).

The test set contains 86 randomly selected voicemail messages (approximately 5 7000 words). For every test message, a firstpass speaker independent decoding produced a MAP word lattice described in section 3. For the MLLR statistics we used phone and gaussian posteriors as described in section 4. The regression classes for MLLR were defined in the following way: first all the mixture components within a phone were bottom-up clustered using a minimum likelihood distance and next, the representatives for 10 all the phones were clustered again until reaching one root node. The number of MLLR transforms that will be computed depends on the number of counts that particular nodes in the regression tree get. In practice, a minimum threshold of 1500 was found to be useful. For voicemail messages which are typically 10 to 50 seconds long this results in computing 13 transforms per message.

15 Figure 3 shows the histogram of the non zero phone posteriors computed over all the test sentences. It is to be noted that, first, there are a significant number of entries with moderate (0.1 - 0.9) probabilities. Secondly, although there are a significant number of entries at the leftend of the histogram, they have such low probabilities that they account

for an insignificant amount of probability mass. This suggests that one can use high values for the confidence thresholds on the posteriors without losing too much adaptation data.

Figure 4 shows the word error rate as a function of the confidence threshold. The optimal results were obtained for a threshold of 0.8. Increasing the threshold above this value results in discarding too much adaptation data which counters the effect of using only alignments in which one is very confident.

Finally, Table I (herebelow) compares the word error rates of the speaker independent system, 1-best MLLR, lattice MLLR and confidence-based lattice MLLR. The overall improvement of the confidence-based lattice MLLR over the 1best MLLR is about 1.8% relative and has been found to be consistent across different test sets. It is expected that the application of iterative MLLR, i.e. repeated data alignment and transform estimation, will increase the improvement. This is because the lattice has more correct words to align to than the 1best transcription. For comparison, Wallhoff et al., "Frame-Discriminative and Confidence-Driven Adaptation for LVCSR" (Proceedings of ICASSP '00, 2000) cites a gain on a "Wall Street Journal" task of 34% relative over standard MLLR by combining confidence measures with MLLR.

TABLE 1

System	Word Error Rate
Baseline (SI)	33.72%
1-best MLLR	32.14%
Lattice MLLR	31.98%
Lattice MLLR + threshold	31.56%

In recapitulation, the present invention, in accordance with at least one presently
5 preferred embodiment, broadly contemplates the use of a word lattice in conjunction with
MLLR. Rather than adjusting the gaussian means to maximize the likelihood of the data
given a single decoded script, a transform was generated that maximized the likelihood of
the data given a set of word hypotheses concisely represented in a word lattice. It was
found that the use of a lattice alone produces an improvement, and also that one can gain a
10 more significant improvement by discarding statistics in which one has low confidence.

In further recapitulation, it will be appreciated from the foregoing that the use of lattice-based information for unsupervised speaker adaptation is explored herein. It is recognized that, as initially formulated, MLLR aims to linearly transform the means of the gaussian models in order to maximize the likelihood of the adaptation data given the

5 correct hypothesis (supervised MLLR) or the decoded hypothesis (unsupervised MLLR). For the latter, if the first-pass decoded hypothesis is significantly erroneous (as is usually the case for large vocabulary telephony applications), MLLR will often find a transform that increases the likelihood for the incorrect models, and may even lower the likelihood of the correct hypothesis. Since the oracle word error rate of a lattice is much lower than

10 that of the 1-best or N-best hypothesis, by performing adaptation against a word lattice, correct models are more likely to be used in estimating a transform. Further, a particular type of lattice proposed herein enables the use of a natural confidence measure given by the posterior occupancy probability of a state, that is, the statistics of a particular state will be updated with the current frame only if the *a posteriori* probability of the state at that

15 particular time is greater than a predetermined threshold. Experiments performed on a voicemail speech recognition task indicate a relative 2% improvement in the word error rate of lattice MLLR over 1-best MLLR.

The present invention is applicable to all particular forms of MLLR, including those in which the gaussian variances are transformed, and those in which the feature vectors are transformed.

It is to be understood that the present invention, in accordance with at least one
5 presently preferred embodiment, includes at least one speech recognition model, an
accepting arrangement which accepts speaker data, a lattice generator which generates a
word lattice based on the speaker data, and a processing arrangement which adapts at
least one of the speaker data and the at least one speech recognition model in a manner to
maximize the likelihood of the speaker data with respect to the generated word lattice.
10 Together, the accepting arrangement, lattice generator and processing arrangement may
be implemented on at least one general-purpose computer running suitable software
programs. These may also be implemented on at least one Integrated Circuit or part of at
least one Integrated Circuit. Thus, it is to be understood that the invention may be
implemented in hardware, software, or a combination of both.

15 If not otherwise stated herein, it is to be assumed that all patents, patent
applications, patent publications and other publications (including web-based publications)
mentioned and cited herein are hereby fully incorporated by reference herein as if set forth
in their entirety herein.

Claims

What is claimed is:

1. A method of providing speaker adaptation in speech recognition, said method comprising the steps of:

5 providing at least one speech recognition model;

 accepting speaker data;

 generating a word lattice based on the speaker data; and

 adapting at least one of the speaker data and the at least one speech recognition
model in a manner to maximize the likelihood of the speaker data with respect to the
10 generated word lattice.

2. The method according to Claim 1, wherein said step of generating a word lattice comprises generating a maximum a-posteriori probability word lattice.

3. The method according to Claim 2, wherein said step of generating a maximum a-posteriori probability word lattice comprises:

determining posterior state occupancy probabilities for each state in the speaker data at each time;

determining posterior word occupancy probabilities by summing over all states interior to each word in the speaker data; and

5 determining at least one likeliest word at each frame of the speaker data.

4. The method according to Claim 2, wherein said step of generating a word lattice further comprises connecting word traces into a lattice.

5. The method according to Claim 1, further comprising the step of discarding interpretations associated with low confidence.

10 6. The method according to Claim 5, wherein said discarding step comprises determining posterior phone probability.

7. The method according to Claim 1, wherein said adapting step comprises performing maximum likelihood linear regression on the speaker data.

8. An apparatus for providing speaker adaptation in speech recognition, said
15 apparatus comprising:

at least one speech recognition model;

an accepting arrangement which accepts speaker data;

a lattice generator which generates a word lattice based on the speaker data; and

a processing arrangement which adapts at least one of the speaker data and the at
5 least one speech recognition model in a manner to maximize the likelihood of the speaker
data with respect to the generated word lattice.

9. The apparatus according to Claim 8, wherein said generator is adapted to
generate a maximum a-posteriori probability word lattice.

10. The apparatus according to Claim 9, wherein said generator is adapted to:
10 determine posterior state occupancy probabilities for each state in the speaker data
at each time;

determine posterior word occupancy probabilities by summing over all states
interior to each word in the speaker data; and

determine at least one likeliest word at each frame of the speaker data.

11. The apparatus according to Claim 9, wherein said generator is further adapted to connect word traces into a lattice.

12. The apparatus according to Claim 8, further comprising a discarding arrangement which discards interpretations associated with low confidence.

5 13. The apparatus according to Claim 12, wherein said discarding arrangement is adapted to determine posterior phone probability.

14. The apparatus according to Claim 8, wherein said processing arrangement is adapted to perform maximum likelihood linear regression on the speaker data.

15. A program storage device readable by machine, tangibly embodying a program
10 of instructions executable by the machine to perform method steps for providing speaker adaptation in speech recognition, said method comprising the steps of:

providing at least one speech recognition model;

accepting speaker data;

generating a word lattice based on the speaker data; and

adapting at least one of the speaker data and the at least one speech recognition model in a manner to maximize the likelihood of the speaker data with respect to the generated word lattice.

LATTICE-BASED UNSUPERVISED MAXIMUM LIKELIHOOD LINEAR
REGRESSION FOR SPEAKER ADAPTATION

Abstract of the Disclosure

Methods and arrangements using lattice-based information for unsupervised
5 speaker adaptation. By performing adaptation against a word lattice, correct models are
more likely to be used in estimating a transform. Further, a particular type of lattice
proposed herein enables the use of a natural confidence measure given by the posterior
occupancy probability of a state, that is, the statistics of a particular state will be updated
with the current frame only if the *a posteriori* probability of the state at that particular time
10 is greater than a predetermined threshold.

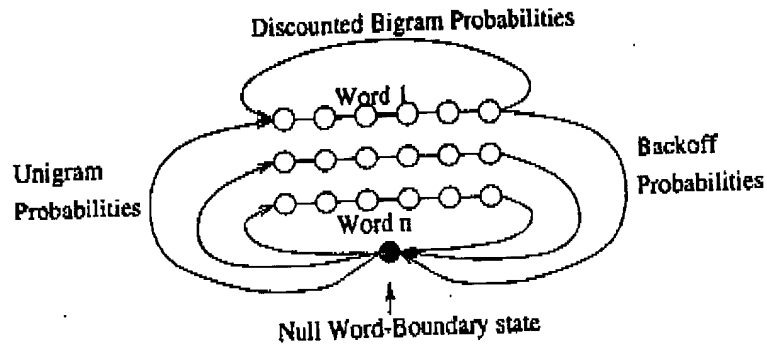


FIG. 1

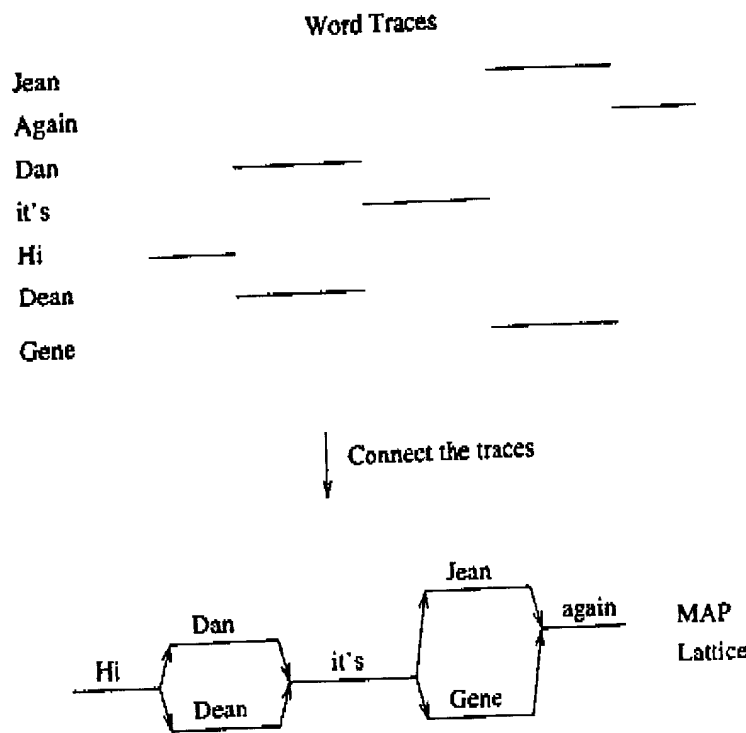


FIG. 2

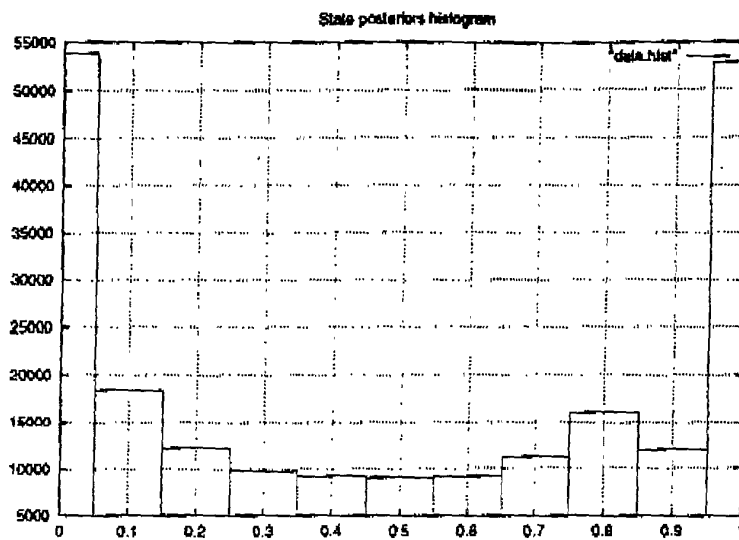


FIG. 3

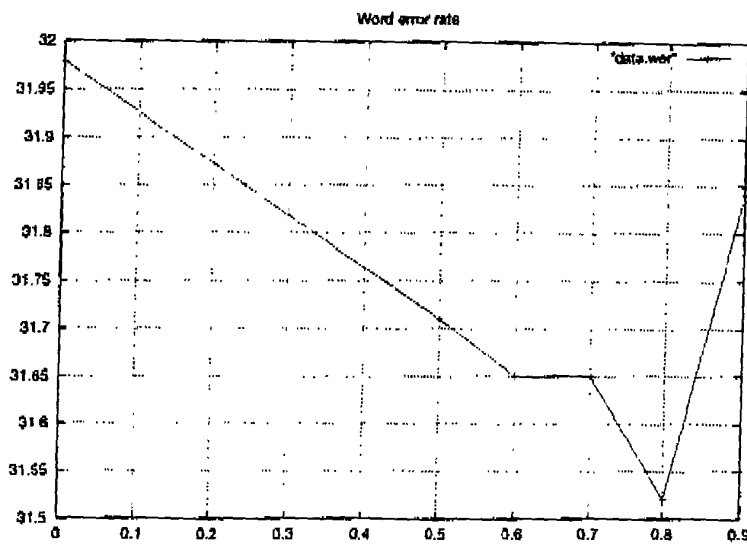


FIG. 4

DECLARATION AND POWER OF ATTORNEY FOR PATENT APPLICATION

As a below named inventor, I hereby declare that:

My residence, post office address and citizenship are as stated below next to my name;

I believe I am the original, first and sole inventor (if only one name is listed below) or an original, first and joint inventor (if plural names are listed below) of the subject matter which is claimed and for which a patent is sought on the invention entitled:

LATTICE-BASED UNSUPERVISED MAXIMUM LIKELIHOOD LINEAR REGRESSION FOR SPEAKER ADAPTATION

the specification of which (check one)

☒ is attached hereto.

_____ was filed on _____ as International Business Machines Docket No. YOR920000390US1

and was amended on _____ (if applicable)

I hereby state that I have reviewed and understand the contents of the above identified specification, including the claims, as amended by any amendment referred to above.

I acknowledge the duty to disclose information which is material to the patentability of this application in accordance with Title 37, Code of Federal Regulations, Section 1.56.

I hereby claim foreign priority benefits under Title 35, United States Code, §119(a)-(d) or §365(b) of any foreign application(s) for patent or inventor's certificate, or §365(a) of any PCT International application which designated at least one country other than the United States, listed below and have also identified below, by checking the box, any foreign application for patent or inventor's certificate, or PCT International application, having a filing date before that of the application on which priority is claimed:

Prior Foreign Application(s)			Priority Claimed	
_____ (Number)	_____ (Country)	_____ (Day/Month/Year Filed)	_____ Yes	_____ No
_____ (Number)	_____ (Country)	_____ (Day/Month/Year Filed)	_____ Yes	_____ No
_____ (Number)	_____ (Country)	_____ (Day/Month/Year Filed)	_____ Yes	_____ No

I hereby claim the benefit under 35 U.S.C. §119(e) of any United States provisional application(s) listed below.

_____ (Application Number)	_____ (Filing Date)
_____ (Application Number)	_____ (Filing Date)

DECLARATION AND POWER OF ATTORNEY FOR PATENT APPLICATION

I hereby claim the benefit under 35 U.S.C. §120 of any United States Application(s), or §365(c) of any PCT International application designating the United States, listed below and, insofar as the subject matter of each of the claims of this application is not disclosed in the prior United States, or PCT International application in the manner provided by the first paragraph of 35 U.S.C. §112, I acknowledge the duty to disclose information material to the patentability of this application as defined in 37 CFR §1.56 which occurred between the filing date of the prior application and the national or PCT international filing date of this application:

_____ (Application Serial No.)	_____ (Filing Date)	_____ (Status) (patented, pending, abandoned)
-----------------------------------	------------------------	--

_____ (Application Serial No.)	_____ (Filing Date)	_____ (Status) (patented, pending, abandoned)
-----------------------------------	------------------------	--

I hereby declare that all statements made herein of my own knowledge are true and that all statements made on information and belief are believed to be true; and further that these statements were made with the knowledge that willful false statements and the like so made are punishable by fine or imprisonment, or both, under Section 1001 of Title 18 of the United States Code and that willful false statements may jeopardize the validity of the application or any patent issued thereon.

POWER OF ATTORNEY: As a named inventor I hereby appoint the following attorney(s) and/or agent(s) to prosecute this application and transact all business in the Patent and Trademark Office connected therewith (list name and registration number).

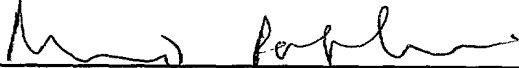
Manny W. Schecter (Reg. 31,722), Terry J. Ilardi (Reg. 29,936), Christopher A. Hughes (Reg. 26,914), Edward A. Pennington (Reg. 32,588), John E. Hoel (Reg. 26,279), Joseph C. Redmond, Jr. (Reg. 18,753), Paul J. Otterstedt (Reg. 37,411), Douglas W. Cameron (Reg. 31,596), Kevin M. Jordan (Reg. 40,277), Stephen C. Kaufman (Reg. 29,551), Daniel P. Morris (Reg. 32,053), Louis J. Percello (Reg. 33,206), Jay P. Sbröllini (Reg. 36,266), David M. Shofi (Reg. 39,835), Robert M. Trepp (Reg. 25,933), and Louis P. Herzberg (Reg. 41,500)

Send Correspondence to: FERENCE & ASSOCIATES, 129 Oakhurst Road, Pittsburgh, PA 15215

Direct Telephone Calls to: (name and telephone number) Stanley D. Ference III, (412) 781-7386

Mukund Padmanabhan

Full name of sole or first inventor



Inventor's Signature

Sep 22, 2000

Date

19 Old Mamaroneck Road, White Plains, NY 10605
Residence

Indian

Citizenship

Same as above

Post Office Address

DECLARATION AND POWER OF ATTORNEY FOR PATENT APPLICATION

George A. Saon

Full name of second joint-inventor, if any

George A. Saon
Inventor's Signature

9/22/2000
Date

142 Kramers Pond Road, Putnam Valley, NY 10579

Residence

Romanian

Citizenship

Same as above

Post Office Address

Geoffrey G. Zweig

Full name of third joint-inventor, if any

Geoffrey G. Zweig
Inventor's Signature

9/22/2000
Date

30 Brookside Drive, Apt. 1A, Greenwich, CT 06830

Residence

USA

Citizenship

Same as above

Post Office Address